

Feasible Automata for Two-Variable Logic with Successor on Data Words^{*}

Ahmet Kara¹, Thomas Schwentick¹, and Tony Tan²

¹ Technical University of Dortmund

² University of Edinburgh

Abstract. We introduce an automata model for data words, that is words that carry at each position a symbol from a finite alphabet and a value from an unbounded data domain. The model is (semantically) a restriction of data automata, introduced by Bojanczyk, et. al. in 2006, therefore it is called *weak data automata*. It is strictly less expressive than data automata and the expressive power is incomparable with register automata. The expressive power of weak data automata corresponds exactly to existential monadic second order logic with successor $+1$ and data value equality \sim , $\text{EMSO}^2(+1, \sim)$. It follows from previous work, David, et. al. in 2010, that the nonemptiness problem for weak data automata can be decided in 2-NEXPTIME. Furthermore, we study weak Büchi automata on data ω -strings. They can be characterized by the extension of $\text{EMSO}^2(+1, \sim)$ with existential quantifiers for infinite sets. Finally, the same complexity bound for its nonemptiness problem is established by a nondeterministic polynomial time reduction to the nonemptiness problem of weak data automata.

1 Introduction

Motivated by challenges in XML reasoning and infinite-state Model Checking, an extension of strings and finitely labelled trees by data values has been investigated in recent years. In classical automata theory, a string is a sequence of positions that carry a symbol from some finite alphabet. In a nutshell, *data strings* generalize strings, in that every position additionally carries a data value from some infinite domain. In the same way, *data trees* generalize (finitely) labelled trees. In XML Theory, data trees model XML documents. Here, the data values can be used to represent attribute values or text content. Both, cannot be adequately modelled by a finite alphabet. In a Model Checking³ scenario, the data values can be used, e.g., to represent process id's or other data.

Early investigations in this area usually considered strings over an “infinite alphabet”, that is, each position only have a value, but no finite-alphabet symbol [2,19,7,14,15,17]. Many of the automata models and logics that have been studied for data strings and trees

^{*} The first and the second authors acknowledge the financial support by the German DFG under grant SCHW 678/4-1, and the third author Future and Emerging Technologies (FET) programme within the Seventh Framework Programme for Research of the European Commission, under the FET-Open grant agreement FOX, number FP7-ICT-233599.

³ In the Model Checking setting, a position might carry a finite set of propositional variables, instead of a symbol.

lack the usual nice decidability properties of automata over finite alphabets, unless strong restrictions are imposed [10,4,3,1].

A result that is particularly interesting for our investigations is the decidability of the satisfiability problem for two-variable logic over data strings [4]. Here, as usual, the logical quantifiers range over the positions of the data string and it can be checked whether a position x carries a symbol a (written: $a(x)$), whether it is to the left of a position y ($x + 1 = y$), whether x is somewhere to the left of y ($x < y$) and whether x and y carry the same data value ($x \sim y$). The logic is denoted by $\text{FO}^2(+1, <, \sim)$. The result was shown with the help of a newly introduced automata model for data words, *data automata* (DA). It turned out, that the expressive power of these automata can be actually characterized by the extension of $\text{FO}^2(+1, <, \sim)$ with existential quantification over sets (of positions) and an additional predicate that holds for x and y if y is the next position from x with the same data value.

However, the complexity of the decision procedure for $\text{FO}^2(+1, <, \sim)$ is very high. The problem is equivalent to the Reachability problem for Petri nets [12], a notoriously hard problem whose complexity has not been resolved exactly. Thus, it has been investigated how the complexity can be reduced, by dropping one of the predicates $x < y$ or $x + 1 = y$. In the latter case (that is, for $\text{FO}^2(<, \sim)$) the complexity decreases to NEXPTIME [4]. In the former case ($\text{FO}^2(+1, \sim)$) the complexity also becomes elementary. In [3] a 3-NEXPTIME bound was shown for the case of data trees and this bound clearly carries over to data strings. A more direct proof with a 4-NEXPTIME bound was given in [8] and a 2-NEXPTIME bound was obtained in [18].

The high complexity of the satisfiability of $\text{FO}^2(+1, <, \sim)$ in [4] results from the high complexity of the nonemptiness problem for data automata. One of the starting questions for this paper was:

- (1) Is there a natural restriction of data automata with (i) a better complexity and (ii) a correspondence to $\text{FO}^2(+1, \sim)$?

We show that such a restriction indeed exists. Data automata consist of two automata \mathcal{A} and \mathcal{B} . \mathcal{A} is a non-deterministic letter-to-letter transducer that constructs, given the finite alphabet part of the input data string⁴ u , a new data string w (where, for each position, the data value in w is the same as in u). The second automaton \mathcal{B} can then check properties of the subsequences of w that carry the same data value. We define *weak data automata* (WDA) which also use a non-deterministic letter-to-letter transducer but can only test some simple constraints of the subsequences in the second part. These constraints are (unary) key, inclusion and denial constraints and they are evaluated for each class separately (there are no inter-class constraints).

It turns out that WDA are expressively weaker than data automata, incomparable with register automata [14,1] and that their expressiveness can be precisely characterized by the extension of $\text{FO}^2(+1, \sim)$ by existential set quantification, that is, $\text{EMSO}^2(+1, \sim)$. As the property that we use to separate the expressive power of WDA and DA can be defined in $\text{EMSO}^2(+1, <, \sim)$ we get that $\text{EMSO}^2(+1, \sim) \not\equiv \text{EMSO}^2(+1, <, \sim)$ as opposed to the classical setting (without data values) where $\text{EMSO}^2(+1) \equiv \text{EMSO}^2(+1, <)$. Indeed, one of the benefits of the logical characterization is that it gives an easy means to show

⁴ The transducer also sees whether a position has the same data value as the next one.

non-expressibility for $\text{EMSO}^2(+1, \sim)$ (and $\text{FO}^2(+1, \sim)$). From results in [8] it immediately follows that the nonemptiness problem for WDA can be solved in 2-NEXPTIME.

As mentioned above, one motivation to study data strings comes from Model Checking. In that context, systems are usually considered to run forever and to produce infinite traces. Thus, data ω -words need to be considered as well, and this was actually one of the main motivations of this research. In particular we address the following questions.

- (2) Do the complexity results of [8] carry over to data ω -strings?
- (3) Can the expressibility results and logical characterizations of the first part of the paper also be established for data ω -strings?

It is straightforward to adapt weak data automata for data ω -strings. The transducer can simply be equipped with a Büchi acceptance mechanism. We refer to the resulting model as *weak Büchi data automata* (WBDA). It turns out that the answer to both questions, (2) and (3), is affirmative. For (3), this is not hard to prove. The separation of WDA from DA also separates WBDA from Büchi data automata. It is also not too hard to get a logical characterization of WBDA by extending $\text{EMSO}^2(+1, \sim)$ with existential set quantifiers that are semantically restricted to bind to infinite sets. The answer to question (2) required considerably more effort. However, we establish a 2-NEXPTIME upper bound for the nonemptiness problem for WBDA by a nondeterministic polynomial time reduction to the nonemptiness for WDA.

Related work. Some related work was already mentioned above. The pioneering works in Linear Temporal Logic for ω -words with data are the papers [10,9]. In [9] an extension of Linear Temporal Logic (LTL) to handle data values is proposed and its satisfiability problem is shown to be decidable. The decision procedure is a reduction to the reachability problem in Petri nets, thus resulting in a similarly unknown complexity as for data automata. The logic and automata considered in [10] are decidable for finite data words, but not primitive recursive, and undecidable for ω -words. In [16] it is shown that with a *safety* restriction both the logic and the automata become decidable, even in EXPSpace. In [9] a logic with PSPACE complexity is considered. In [5], MSO logic on data words (with possibly multiple data values per position) is compared to automata models for various types of successor relations.

Organization. We give basic definitions in Section 2. In Section 3, weak data automata are defined, their complexity is given, and their expressive power is compared with other models. Section 4 gives the logical characterization of WDA by $\text{EMSO}^2(+1, \sim)$. Section 5 studies data ω -strings and shows how the nonemptiness problem of WBDA can be nondeterministically reduced in polynomial time to the nonemptiness of WDA. Section 6 states some open problems. .

Acknowledgement. We thank Christof Löding for helpful remarks on automata and logics for ω -words and Thomas Zeume for thorough proof reading.

2 Notation

Data words. Let Σ be a finite alphabet and \mathcal{D} an infinite set of data values. A *finite* word is an element of Σ^* , while an ω -word is an element of Σ^ω . A finite *data word* is an element

of $(\Sigma \times \mathfrak{D})^*$, while a *data ω -word* is an element of $(\Sigma \times \mathfrak{D})^\omega$. We often refer to data words also as *data strings*.

We write a data (finite or ω -) word w as $\binom{a_1}{d_1} \binom{a_2}{d_2} \cdots$, where $a_1, a_2, \dots \in \Sigma$ and $d_1, d_2, \dots \in \mathfrak{D}$. The symbol a_i is the label of position i , while the value d_i is the data value of position i . The projection of w to the alphabet Σ is denoted by $\text{Str}(w) = a_1 a_2 \dots$. A position in w is called an a -position, if the label of that position is a . We denote by $V_w(a)$, the set of data values found in a -positions in w , i.e., $V_w(a) = \{d_i \mid a_i = a\}$, for each $a \in \Sigma$. Note that some $V_w(a)$'s may be infinite, while some others finite.

A maximal set of positions with the same data value d is called a *class* c^d of the word and the Σ -string induced by the symbols at its positions is called the *class string* w^d . The *profile word* of a data ω -word $w = \binom{a_1}{d_1} \binom{a_2}{d_2} \cdots$ is $\text{Profile}(w) = (a_1, s_1), (a_2, s_2), \dots \in (\Sigma \times \{\top, \perp\})^\omega$, where for each position $i \geq 1$ the component s_i is \top if and only if $d_i = d_{i+1}$. The profile word of a finite data word $\binom{a_1}{d_1} \binom{a_2}{d_2} \cdots \binom{a_n}{d_n}$ is defined similarly, with the addition that the component s_n is \perp .

Automata and Büchi automata. An *automaton* \mathcal{A} over the alphabet Σ is a tuple $\mathcal{A} = \langle \Sigma, Q, q_0, \Delta, F \rangle$, where Q is a finite set of states, $q_0 \in Q$ is the initial state, $\Delta \subseteq Q \times \Sigma \times Q$ is a set of transitions and $F \subseteq Q$ is a set of accepting states. A run of \mathcal{A} on a word $w = a_1 a_2 \dots a_n$ is a sequence $\rho = q_1 \dots q_n$ of states from $Q - \{q_0\}$ such that $(q_0, a_1, q_1) \in \Delta$ and $(q_i, a_{i+1}, q_{i+1}) \in \Delta$ for each $1 \leq i < n$. The run ρ is accepting, if $q_n \in F$.

A *Büchi automaton* \mathcal{A} is syntactically just an automaton. A run of \mathcal{A} on an ω -word $w = a_1 a_2 \dots$ is an infinite sequence $\rho = q_1 q_2 \dots$ of states from $Q - \{q_0\}$ such that $(q_0, a_1, q_1) \in \Delta$ and $(q_i, a_{i+1}, q_{i+1}) \in \Delta$, for each $i \geq 1$. Let $\text{Inf}(\rho)$ denote the set of states that appear infinitely many times in ρ . The run ρ is accepting if $\text{Inf}(\rho) \cap F \neq \emptyset$.

A word (resp. an ω -word) w is accepted by an automaton (resp. Büchi automaton) \mathcal{A} , if there exists an accepting run of \mathcal{A} on w . As usual, $\mathcal{L}(\mathcal{A})$ (resp. $\mathcal{L}^\omega(\mathcal{A})$) denotes the set of words (resp. ω -words) accepted by the automaton \mathcal{A} .

Letter-to-letter transducers. A *letter-to-letter transducer* over the input alphabet Σ and the output alphabet Γ is a tuple $\mathcal{T} = \langle \Sigma, \Gamma, Q, q_0, \Delta, F \rangle$, where Q, q_0, F are the set of states, the initial state, and the set of accepting states, respectively, and $\Delta \subseteq Q \times \Sigma \times Q \times \Gamma$ is the set of transitions. The intuitive meaning of a transition (q, a, q', γ) is that when the automaton is in state q , reading the symbol a , then it can move to the state q' and output γ . A *run* of \mathcal{T} on a word $w = a_1 a_2 \dots a_n$ is a sequence $(q_1, \gamma_1), \dots, (q_n, \gamma_n)$ over $(Q - \{q_0\}) \times \Gamma$ such that $(q_0, a_1, q_1, \gamma_1) \in \Delta$ and $(q_i, a_{i+1}, q_{i+1}, \gamma_{i+1}) \in \Delta$, for each $1 \leq i < n$. Likewise, a *run* of \mathcal{T} on an ω -word $w = a_1 a_2 \dots$ is a sequence $(q_1, \gamma_1), (q_2, \gamma_2), \dots$ over $(Q - \{q_0\}) \times \Gamma$ such that $(q_0, a_1, q_1, \gamma_1) \in \Delta$ and $(q_i, a_{i+1}, q_{i+1}, \gamma_{i+1}) \in \Delta$, for each $i \geq 1$. A run is *accepting* if it is accepting in the sense of (Büchi) automata. We say that $v = \gamma_1 \gamma_2 \dots$ is an output of \mathcal{T} on w , if there exists an accepting run $(q_1, \gamma_1), (q_2, \gamma_2), \dots$ of \mathcal{T} on w .

Data automata. A *data automaton (DA)* is a pair $(\mathcal{A}, \mathcal{B})$, where

- \mathcal{A} is a letter-to-letter transducer with input alphabet $\Sigma \times \{\top, \perp\}$ and output alphabet Γ ,
- \mathcal{B} is a finite state automaton over the alphabet Γ .

A data word w is accepted by $(\mathcal{A}, \mathcal{B})$ if the following holds.

- $\text{Profile}(w)$ is accepted by \mathcal{A} , yielding an output u .
- For each data value d of w , the class string u^d is accepted by \mathcal{B} .

Data automata were introduced in the stated form in [4]. In [1] it was shown that their expressive power is not affected, if \mathcal{A} gets $\text{Str}(w)$ as input as opposed to $\text{Profile}(w)$. In more recent papers, data automata are therefore defined in the (syntactically) weaker form with input $\text{Str}(w)$.

3 Weak data automata

In this section we define a new automata model for finite data words and study its expressive power and its complexity. The model follows a similar approach as the model of data automata. The profile of the input data word is transformed by a letter-to-letter transducer and then further conditions on the resulting class strings are imposed. However, the conditions that can be stated in the new automata model are much more limited than those of a data automaton (hence the name *weak* data automata).

Let Γ be an alphabet. Weak data automata allow three kinds of data constraints over Γ :

1. *key constraints*, written in the form: $\text{key}(\gamma)$, where $\gamma \in \Gamma$.
2. *inclusion constraints*, written in the form: $V(\gamma) \subseteq \bigcup_{\gamma' \in R} V(\gamma')$, where $\gamma \in \Gamma$, $R \subseteq \Gamma$.
3. *denial constraints*, written in the form: $V(\gamma) \cap V(\gamma') = \emptyset$, where $\gamma, \gamma' \in \Gamma$.

Whether a data word w satisfies a data constraint C , written as $w \models C$, is defined as follows.

1. $w \models \text{key}(\gamma)$, if every two γ -positions in w have different data values.
2. $w \models V(\gamma) \subseteq \bigcup_{\gamma' \in R} V(\gamma')$, if $V_w(\gamma) \subseteq \bigcup_{\gamma' \in R} V_w(\gamma')$.
3. $w \models V(\gamma) \cap V(\gamma') = \emptyset$, if $V_w(\gamma) \cap V_w(\gamma') = \emptyset$.

If \mathcal{C} is a collection of data constraints, then we write $w \models \mathcal{C}$, if $w \models C$ for all $C \in \mathcal{C}$.

A *weak data automaton (WDA)* over the alphabet Σ is a pair $(\mathcal{A}, \mathcal{C})$, where

- \mathcal{A} is a letter-to-letter transducer with input alphabet $\Sigma \times \{\top, \perp\}$ and output alphabet Γ ,
- \mathcal{C} is a collection of data constraints over the alphabet Γ .

A data word $w = (a_1)_{d_1} (a_2)_{d_2} \dots (a_n)_{d_n}$ is accepted by a WDA $(\mathcal{A}, \mathcal{C})$, if

- there is an accepting run of \mathcal{A} on $\text{Profile}(w)$, with an output $\gamma_1 \gamma_2 \dots \gamma_n$, and
- the induced data word $w = (\gamma_1)_{d_1} (\gamma_2)_{d_2} \dots (\gamma_n)_{d_n}$ satisfies all the constraints in \mathcal{C} .

We write $\mathcal{L}(\mathcal{A}, \mathcal{C})$ to denote the language that consists of all data words accepted by $(\mathcal{A}, \mathcal{C})$.

We first discuss some extensions of WDA by the constraints that were studied in [8].

- *Disjunctive key constraints* are written in the form: $\text{key}(K)$, where $K \subseteq \Gamma$. Such a constraint is satisfied by a data word if each of its classes has at most one position with a symbol from K .

- *Disjunctive inclusion constraints* are written in the form: $\bigcup_{\gamma \in S} V(\gamma) \subseteq \bigcup_{\gamma' \in R} V(\gamma')$, where $S, R \subseteq \Gamma$. Such a constraint is satisfied by a data word if each class with a position with a symbol from S also has a position with a symbol from R .

An *extended weak data automaton* is defined like a WDA but it further allows disjunctive key and inclusion constraints.

Lemma 1. *From each extended WDA $(\mathcal{A}, \mathcal{C})$ an equivalent WDA of polynomial size can be constructed in polynomial time.*

Proof. A disjunctive inclusion constraint $\bigcup_{\gamma \in S} V(\gamma) \subseteq \bigcup_{\gamma' \in R} V(\gamma')$ can simply be replaced by a set of inclusion constraints $V(\gamma) \subseteq \bigcup_{\gamma' \in R} V(\gamma')$, one for each $\gamma \in S$. Disjunctive key constraints $\text{key}(K)$ can be replaced by a set of denial constraints $V(\gamma) \cap V(\gamma') = \emptyset$, one for each pair $\gamma \neq \gamma'$ with $\gamma, \gamma' \in K$ and a set of key constraints $\text{key}(\gamma)$, one for each $\gamma \in K$. \square

Next, we compare the expressive power of weak data automata with other automata models for data words. More precisely we compare it with register automata [14,1] and data automata. Register automata are an extension of finite state automata with a fixed number of registers in which they can store data values and compare them with the data value of subsequent positions. For a precise definition we refer¹ the reader to [1].

We consider the following two data languages.

- $L_{a < b}$ consists of all data words over the alphabet $\{a, b\}$ with the property that for every a -position i there is a b -position $j > i$ with the same data value;
- $L_{a * b}$ is the subset of $L_{a < b}$ where the next b -position j with the same data value as i always satisfies $j = i + 2$.

Lemma 2. *Neither $L_{a * b}$ nor $L_{a < b}$ can be decided by a WDA.*

Proof. We first show that no WDA decides $L_{a * b}$. Towards a contradiction, we thus assume that $L_{a * b}$ is decided by some weak data automata $(\mathcal{A}, \mathcal{C})$.

To this end, let $n = |\Gamma|^4 + 1$ and let $d_1, d'_1, d_2, d'_2, \dots, d_n, d'_n$ be pairwise different data values. We consider the data word

$$w = \begin{pmatrix} a \\ d_1 \end{pmatrix} \begin{pmatrix} a \\ d'_1 \end{pmatrix} \begin{pmatrix} b \\ d_1 \end{pmatrix} \begin{pmatrix} b \\ d'_1 \end{pmatrix} \begin{pmatrix} a \\ d_2 \end{pmatrix} \begin{pmatrix} a \\ d'_2 \end{pmatrix} \begin{pmatrix} b \\ d_2 \end{pmatrix} \begin{pmatrix} b \\ d'_2 \end{pmatrix} \cdots \begin{pmatrix} a \\ d_n \end{pmatrix} \begin{pmatrix} a \\ d'_n \end{pmatrix} \begin{pmatrix} b \\ d_n \end{pmatrix} \begin{pmatrix} b \\ d'_n \end{pmatrix}$$

of length $4n$. Clearly, w is in $L_{a * b}$ and its profile is $((a, \perp)(a, \perp)(b, \perp)(b, \perp))^n$.

Let $\gamma = \gamma_1 \gamma_2 \cdots \gamma_{4n}$ be an output of \mathcal{A} on the profile of w such that $\begin{pmatrix} \gamma_1 \\ d_1 \end{pmatrix} \cdots \begin{pmatrix} \gamma_{4n} \\ d'_n \end{pmatrix}$ satisfies all constraints in \mathcal{C} . By the choice of n , there exist numbers i, j with $0 \leq i < j < n$ such that $\gamma_{4i+1} \gamma_{4i+2} \gamma_{4i+3} \gamma_{4i+4} = \gamma_{4j+1} \gamma_{4j+2} \gamma_{4j+3} \gamma_{4j+4}$.

Let u be the data word obtained from w by swapping the positions of the data values $d_{i+1} d'_{i+1}$ and $d_{j+1} d'_{j+1}$. That is,

$$u = \begin{pmatrix} a \\ d_1 \end{pmatrix} \cdots \begin{pmatrix} a \\ d_{i+1} \end{pmatrix} \begin{pmatrix} a \\ d'_{i+1} \end{pmatrix} \begin{pmatrix} b \\ d_{j+1} \end{pmatrix} \begin{pmatrix} b \\ d'_{j+1} \end{pmatrix} \cdots \begin{pmatrix} a \\ d_{j+1} \end{pmatrix} \begin{pmatrix} a \\ d'_{j+1} \end{pmatrix} \begin{pmatrix} b \\ d_{i+1} \end{pmatrix} \begin{pmatrix} b \\ d'_{i+1} \end{pmatrix} \cdots \begin{pmatrix} b \\ d_n \end{pmatrix}.$$

¹ The precursor model *finite-memory automata* was defined on “strings over infinite alphabets”, that is, essentially data strings without a Σ -component [14].

Clearly, $u \notin L_{a*b}$. However, because $\text{Profile}(u) = \text{Profile}(w)$, $\gamma_1\gamma_2 \dots \gamma_{4n}$ is also an output of \mathcal{A} on $\text{Profile}(u)$. Moreover, the sets of $V_u(\gamma) = V_w(\gamma)$, for each $\gamma \in \Gamma$, and therefore the validity of inclusion and denial constraints does not change. Furthermore, as in u and w every data value occurs at exactly one a -position and at exactly one b -position, they cannot be distinguished by key constraints, either. Thus, $u \in \mathcal{L}(\mathcal{A}, \mathcal{C})$, the desired contradiction.

The proof for $L_{a<b}$ is exactly the same, as $w \in L_{a<b}$ and $u \notin L_{a<b}$ (because of $\binom{a}{d_{j+1}}$). \square

Theorem 1. (a) *The class of data languages that are decided by WDA is strictly included in the class of data languages decided by DA.*

(b) *The classes of languages decided by WDA and by register automata are incomparable.*

Proof. Towards (a) we first show that every WDA can be translated into a DA and thus WDA decide a subclass of DA. That the subclass is strict can then be concluded from (b) as register automata are captured by DA [1] and thus there is a data language that can be decided by a DA but not a WDA.

Let thus $(\mathcal{A}, \mathcal{C})$ be a WDA. Then $(\mathcal{A}, \mathcal{B})$ is a data automaton for $L(\mathcal{A}, \mathcal{C})$, where the automaton \mathcal{B} tests the constraints in \mathcal{C} as follows.

- For every key constraint $\text{key}(\gamma)$ of \mathcal{C} , \mathcal{B} tests that every class string has at most one γ -position.
- For every inclusion constraint $V(\gamma) \subseteq \bigcup_{\gamma' \in R} V(\gamma')$, \mathcal{B} tests that every class string with a γ -position also has a γ' -position, for some $\gamma' \in R$.
- For every denial constraint $V(\gamma) \cap V(\gamma') = \emptyset$, \mathcal{B} checks that classes with a γ -position do not have any γ' -positions.

To show statement (b) we first consider the separation language $L = L_{a*b}$ which cannot be decided by a WDA by Lemma 2. However, L_{a*b} can be easily decided by a register automaton that always stores the last two data values in two registers and the information about their symbols in its state.

On the other hand, the set of all data strings over $\Sigma = \{a\}$ in which every data value occurs only once can easily be decided by a WDA by the identity-transducer and the key constraint $\text{key}(a)$ but not by a register automaton [14]. \square

The complexity of the nonemptiness problem for WDA follows directly from results in [8].

Theorem 2. *The nonemptiness problem for WDA is decidable in 2-NEXPTIME.*

Proof. In [8], it was shown that given an automaton \mathcal{A} that reads profile strings and a set \mathcal{C} of disjunctive key and inclusion constraints, to decide whether there is a data word w such that \mathcal{A} accepts $\text{Profile}(w)$ and $w \models \mathcal{C}$ can be done in nondeterministic double exponential time.

Clearly, this is basically the same as the nonemptiness problem for WDA with disjunctive key and inclusion constraints only. It thus only remains to show that denial constraints can be translated into disjunctive constraints in a nonemptiness respecting fashion. To this end, a denial constraint $V(\gamma_1) \cap V(\gamma_2) = \emptyset$ can be replaced as follows. We add two new symbols γ'_1, γ'_2 and require that in each class with γ_i one γ'_i occurs but γ'_1 and γ'_2 do not co-occur by two inclusion constraints $V(\gamma_1) \subseteq V(\gamma'_1)$ and $V(\gamma_2) \subseteq V(\gamma'_2)$ and a disjunctive key constraint for $\{\gamma'_1, \gamma'_2\}$. \square

4 A logical characterization of weak data automata

In this section, we give a logical characterization of the data languages decided by weak data automata in terms of existential second order logic. The characterization is an analogue of the Theorem of Büchi, Elgot and Trakhtenbrot [6,11,21] for string languages. This theorem can be stated for various logics, the most interesting one for our context is that $\text{EMSO}^2(+1)$ characterizes exactly the regular languages.

We represent data words by logical structures $w = \langle \{1, \dots, n\}, +1, <, \{a(\cdot)\}_{a \in \Sigma}, \sim \rangle$, where $\{1, \dots, n\}$ is the set of positions, $+1$ is the successor relation (i.e., $+1(i, j)$ if $i+1 = j$), $<$ is the order relation (i.e., $<(i, j)$ if $i < j$), the $a(\cdot)$'s are the label relations, and $i \sim j$ holds if positions i and j have the same data value. As the empty data word can not be properly represented, the logical characterization of WDA ignores empty data words. That is, if some WDA $(\mathcal{A}, \mathcal{C})$ accepts the empty data string then its language is different from the language of the corresponding formula φ : $\mathcal{L}(\mathcal{A}, \mathcal{C}) = L(\varphi) \cup \{\epsilon\}$.

For a set $SS \subseteq \{+1, <, \sim\}$ of relation symbols, we write $\text{FO}(SS)$ for first-order logic with the vocabulary SS , $\text{MSO}(SS)$ for monadic second-order logic (which extends $\text{FO}(SS)$ with quantification over sets of positions), and $\text{EMSO}(SS)$ for existential monadic second order logic, that is, all sentences of the form $\exists R_1 \dots \exists R_m \psi$, where ψ is an $\text{FO}(SS)$ formula extended with the unary predicates R_1, \dots, R_m . By $\text{FO}^2(SS)$ we denote the restriction of $\text{FO}(SS)$ to sentences with two variables x and y , and by $\text{EMSO}^2(SS)$ the restriction of $\text{EMSO}(SS)$ where the first-order part uses only two variables.

4.1 From weak data automata to $\text{EMSO}^2(+1, \sim)$

Theorem 3. *For every weak data automaton $(\mathcal{A}, \mathcal{C})$, an equivalent $\text{EMSO}^2(+1, \sim)$ -formula φ is constructible in polynomial time.*

Proof. Let $(\mathcal{A}, \mathcal{C})$ be a weak data automaton with $\mathcal{A} = \langle \Sigma, \Gamma, Q, q_0, \Delta, F \rangle$, where $Q = \{q_1, \dots, q_n\}$ and $\Gamma = \{\gamma_1, \dots, \gamma_l\}$. We recall that we assume without loss of generality that \mathcal{A} uses q_0 only its initial state.

We will construct an $\text{EMSO}^2(+1, \sim)$ -formula φ with $L(\mathcal{A}, \mathcal{C}) - \{\epsilon\} = L(\varphi)$. The construction is the same as the classical translation from NFAs to MSO formulas. See, for example, [20].

The formula φ is

$$\varphi = \exists R_{q_1} \dots \exists R_{q_n} \exists R_{\gamma_1} \dots \exists R_{\gamma_l} (\varphi_{\text{part}} \wedge \varphi_{\text{start}} \wedge \varphi_{\text{trans}} \wedge \varphi_{\text{accept}} \wedge \varphi_{\text{constr}})$$

where

- φ_{part} asserts in a straightforward manner that R_{q_1}, \dots, R_{q_n} on one hand and $R_{\gamma_1}, \dots, R_{\gamma_l}$, on the other hand, partition the positions of the input word.
- φ_{start} asserts that the automaton starts in state q_0 :

$$\forall x (\neg \exists y \ x = y + 1 \rightarrow \phi)$$

where ϕ is:

$$\forall y \ y = x + 1 \rightarrow \bigwedge_{a \in \Sigma} \left(\begin{array}{l} (a(x) \wedge x \sim y) \rightarrow \bigvee_{(q_0, (a, \top), q, \gamma) \in \Delta} (R_q(x) \wedge R_\gamma(x)) \\ (a(x) \wedge x \not\sim y) \rightarrow \bigvee_{(q_0, (a, \perp), q, \gamma) \in \Delta} (R_q(x) \wedge R_\gamma(x)) \end{array} \right)$$

– φ_{trans} asserts that transitions are simulated correctly:

$$\forall x \forall y \ y = x + 1 \rightarrow \bigwedge_{a \in \Sigma, q \in Q} \left(\begin{array}{l} (a(y) \wedge R_q(x) \wedge \exists x(x = y + 1 \wedge y \sim x)) \rightarrow \bigvee_{(q, (a, \top), q', \gamma) \in \Delta} R_{q'}(y) \wedge R_\gamma(y) \\ (a(y) \wedge R_q(x) \wedge \exists x(x = y + 1 \wedge y \not\sim x)) \rightarrow \bigvee_{(q, (a, \perp), q', \gamma) \in \Delta} R_{q'}(y) \wedge R_\gamma(y) \end{array} \right)$$

– φ_{accept} states the accepting condition:

$$\forall x [\neg \exists y \ y = x + 1 \rightarrow \bigvee_{q \in F} R_q(x)]$$

– φ_{constr} is a conjunction $\bigwedge_{C \in \mathcal{C}} \psi_C$ over all constraints C from \mathcal{C} such that,

- if C is a key constraint $\text{key}(\gamma)$, then

$$\psi_C = \forall x \forall y [(R_\gamma(x) \wedge R_\gamma(y) \wedge x \sim y) \rightarrow x = y]$$

- if C is an inclusion constraint $V(\gamma) \subseteq \bigcup_{\gamma' \in S} V(\gamma')$, then

$$\psi_C = \forall x \exists y [R_\gamma(x) \rightarrow (\bigvee_{\gamma' \in S} R_{\gamma'}(y) \wedge x \sim y)]$$

- if C is a denial constraint $V(\gamma) \cap V(\gamma') = \emptyset$, then

$$\psi_C = \forall x \forall y [(R_\gamma(x) \wedge R_{\gamma'}(y)) \rightarrow x \not\sim y]$$

The length of φ is $\mathcal{O}(|\Sigma||Q||\Delta| + |\mathcal{C}|)$. The correctness is straightforward and, thus, omitted. \square

4.2 From $\text{EMSO}^2(+1, \sim)$ to weak data automata

In the following, we use the abbreviation $F(x, y)$ for the formula $\neg y = x + 1 \wedge \neg x = y + 1 \wedge x \neq y$, which states that the distance of x and y is at least two.

Theorem 4. *There is an algorithm that translates every $\text{EMSO}^2(+1, \sim)$ -formula φ into an equivalent weak data automaton $(\mathcal{A}, \mathcal{C})$ in doubly exponential time. In particular, the output alphabet Γ of \mathcal{A} and the number of constraints in \mathcal{C} is at most exponential.*

Proof. In the first step, the algorithm transforms φ into an equivalent $\text{EMSO}^2(+1, \sim)$ formula in Scott normal form (SNF) of the form

$$\psi = \exists R_1 \dots \exists R_n [\forall x \forall y \chi' \wedge \bigwedge_{i=1}^m \forall x \exists y \chi'_i],$$

where χ' and each χ'_i are quantifier-free [13]. The size of ψ is linear in the size of φ , in particular, $n = \mathcal{O}(|\varphi|)$ and $m = \mathcal{O}(|\varphi|)$.

Then it rewrites formula χ' into an, at most exponential, conjunction

$$\chi = \bigwedge_j \neg(\alpha_j(x) \wedge \beta_j(y) \wedge \delta_j(x, y) \wedge \epsilon_j(x, y)),$$

where, for every j , α_j, β_j are conjunctions of literals with unary relation symbols, δ_j is either $x \sim y$ or $x \not\sim y$ and $\epsilon_j(x, y)$ is one² of $x = y$, $y = x + 1$, $F(x, y)$.

Likewise, it rewrites every χ'_i into an, at most exponential, disjunction

$$\chi_i = \bigvee_j (\alpha_j^i(x) \wedge \beta_j^i(y) \wedge \delta_j^i(x, y) \wedge \epsilon_j^i(x, y)),$$

where the atomic formulas are of the respective forms as above.

The idea of the construction is that \mathcal{A} guesses³ a couple of relations that allow to state some of the properties expressed in ψ by constraints of \mathcal{C} .

For simplicity we refer to the label relations a_1, \dots, a_l as R_{n+1}, \dots, R_{n+l} .

The relations that are guessed are the following.

- R_1, \dots, R_n (the *SNF relations*). We refer to the full atomic type of a position with respect to the relations R_1, \dots, R_{n+l} as its *SNF-type*;
- P_1, P_2, P_3 with the following intention: if a class contains at least three positions of some SNF-type α , then one of them is in P_3 . If a class contains at least two α -positions, then one of them is in P_2 . If there is at least one α -position then there is one in P_1 .
- C_1, C_2, C_3 with the following intention: if there are at least three classes that contain positions of some SNF-type α , then in one of these classes all α -positions are in C_3 . If there are at least two classes that contain α -positions, then in one of them all α -positions are in C_2 . If there is at least one class with α -positions then in one of them all α -positions are in C_1 . We refer to the full atomic type of a position with respect to $P_1, P_2, P_3, C_1, C_2, C_3$ as its *occurrence type*;
- $E_{\leftarrow}, E_{\rightarrow}$ with the intention that a position is in E_{\leftarrow} if its left neighbor has the same data value (and likewise for E_{\rightarrow});
- $R_1^\ell, \dots, R_{n+l}^\ell, P_1^\ell, P_2^\ell, P_3^\ell, C_1^\ell, C_2^\ell, C_3^\ell$ and $R_1^r, \dots, R_{n+l}^r, P_1^r, P_2^r, P_3^r, C_1^r, C_2^r, C_3^r$ with the following intention: For each position p , it should hold that p is in a relation with superscript ℓ , if its left neighbor is in the corresponding relation without superscript. Likewise, p is in a relation with superscript r , if its right neighbor is in the corresponding relation without superscript. We refer to the type of a position with respect to these relation and the relations $E_{\leftarrow}, E_{\rightarrow}$ as its *neighborhood type*;

² The case $x = y + 1$ does not need to be considered as it can be obtained by swapping x and y .

³ More precisely, \mathcal{A} guesses, for each position p , the set of those relations that contain p . However, on a global level, we refer to this as “guessing the relations”.

- for every $j \leq n+l$ and $i \leq m$ the relations W_j^i and $E^i, G_{suc}^i, G_{pre}^i, G_F^i$ with the following intention: if for some position p , formula χ_i becomes true for $x = p$ and some position $y = q$ then $p \in W_j^i$ if and only if $q \in R_j$. That is, W_1^i, \dots, W_{n+l}^i mimics the SNF-type of witness positions with respect to χ_i . Furthermore, $p \sim q$ if and only if $p \in E^i$, $q = p+1$ if and only if $p \in G_{suc}^i$, $p = q+1$ if and only if $p \in G_{pre}^i$, and $F(p, q)$ if and only if $p \in G_F^i$. We refer to the type with respect to these relations as the *witness type*.

It should be noted that the number of these relations is $\mathcal{O}(nm)$.

Now we describe how \mathcal{A} and \mathcal{C} can be constructed in order to test whether a data string w satisfies φ .

First, the automaton \mathcal{A} guesses which types will occur in the output and verifies during its computation that exactly these types occur. What needs to be verified is how \mathcal{A} and \mathcal{C} can ensure that the relations guessed by \mathcal{A} are consistent with respect to the intention that was described above.

- The consistency with respect to P_1, P_2 and P_3 can be tested as follows. \mathcal{A} can ensure that each position is in at most one P_k . The intention of P_1, P_2 and P_3 can be enforced by the inclusion constraints
 - $V(\alpha \wedge \neg P_1 \wedge \neg P_2) \subseteq V(\alpha \wedge P_3)$,
 - $V(\alpha \wedge P_3) \subseteq V(\alpha \wedge P_2)$, and
 - $V(\alpha \wedge P_2) \subseteq V(\alpha \wedge P_1)$.
 That each class contains at most one position with $\alpha \wedge P_k$, for each k , can be stated by key constraints.
- The consistency with respect to C_1, C_2 and C_3 can be tested in a similar fashion. That, for some α the existence of an $(\alpha \wedge C_3)$ -class implies the existence of an $(\alpha \wedge C_2)$ -class and the other two corresponding conditions can be already guaranteed when \mathcal{A} guesses the set of occurring types. That all α -positions in an $(\alpha \wedge C_k)$ -class are in C_k can be ensured by denial constraints.
- The consistency of the neighborhood types can be easily tested by \mathcal{A} with the help of the profile information.
- For the consistency of the witness types \mathcal{A} checks that for each position p and each $i \leq m$ there actually exists a disjunct $(\alpha_j^i(x) \wedge \beta_j^i(y) \wedge \delta_j^i(x, y) \wedge \epsilon_j^i(x, y))$ of χ_i , for which α_j^i is the SNF-type of p and β_j^i, δ_j^i and ϵ_j^i coincide with the i -th witness type of p . How the existence of corresponding witness positions is tested will be described below.

Now we describe how the conjuncts of χ can be checked. For every conjunct (indexed by j) we distinguish the following cases depending on the possible formulas $\delta_j(x, y)$ and $\epsilon(x, y)$.

- (Case 1) ϵ is $x = y$: in this case the conjunct states a condition about forbidden SNF-types of positions. This kind of constraints can be ensured by \mathcal{A} by not allowing to guess them.
- (Case 2) ϵ is $y = x+1$: such conjuncts state that some pairs of SNF-types are forbidden for neighbors with equal (or different) data values. As this is a question of consistency between the neighborhood type and the witness type of a position it can be guaranteed by \mathcal{A} (by disallowing certain combinations).
- (Case 3) ϵ is $F(x, y)$ and δ is $x \sim y$: such a formula states that there should not be an α -position p and a β -position q in the same class with $|p - q| > 1$. Such a formula

gives rise to some denial constraints. As an example, if the SNF-type of a position p is α and the neighborhood type of p indicates that it has a $(P_1 \wedge \beta)$ -position and a $(P_2 \wedge \beta)$ -position as neighbors then there is a denial constraint forbidding $(P_3 \wedge \beta)$ -positions in this class.

(Case 4) ϵ is $F(x, y)$ and δ is $x \not\sim y$: such a formula states that α -positions p and β -positions q with $|p - q| > 1$ need to be in the same class. This can be tested by \mathcal{A} and some denial constraints with the help of C_1, C_2, C_3 and P_1, P_2, P_3 .

What remains to be shown is how the formulas χ_i can be tested. For each position p and each $i \leq m$, $(\mathcal{A}, \mathcal{C})$ has to check that there is a witness position q such that p and q satisfy some disjunct of χ_i . Which disjunct should be considered is given by the i -th witness type of p . The existence of corresponding positions can be tested in a way that is similar to the tests for formula χ . If ϵ is $x = y$ the witness needs to be p itself which can be tested by \mathcal{A} directly. If ϵ is $y = x + 1$ or $x = y + 1$ the witness is one of the neighbors of p . Whether this is true can be concluded from the neighborhood type and thus also by \mathcal{A} . If ϵ is $F(x, y)$ the existence of witnesses can be checked by inclusion constraints if δ is $x \sim y$ and by \mathcal{A} directly if δ is $x \not\sim y$. In the former case, the P_k -relations are used, in the latter case, both the P_k and the C_k -positions.

The size of the output alphabet of \mathcal{A} and the number of constraints are at most exponential in $|\varphi|$. The number of states is at most doubly exponential. \square

We note that in the upper bound of the algorithm for nonemptiness of WDA transferred from [8], the doubly exponential term only depends on the alphabet size. By combining this with the bounds of Theorem 4 we obtain a 3-NEXPTIME upper bound for satisfiability of $\text{FO}^2(+1, \sim)$ (which is worse than the bound in [18]). We also note that the construction underlying the proof of Theorem 4 can be turned into a nondeterministic exponential time reduction from satisfiability for $\text{FO}^2(+1, \sim)$ to nonemptiness for WDA resulting in an automaton with a singly exponential number of states. The reduction guesses the order in which types appear in the accepted string (as opposed to the construction in the proof of Theorem 4).

The previous two theorems yield the following logical characterization of WDA.

Theorem 5. *Weak data automata and $\text{EMSO}^2(+1, \sim)$ are equivalent in expressive power.*

We note that on strings $\text{EMSO}^2(+1)$ and $\text{EMSO}^2(+1, <)$ are expressively equivalent. It is an interesting consequence of the above characterization that this equivalence does not hold for data strings.

Corollary 1. *The logic $\text{EMSO}^2(+1, \sim)$ is strictly less expressible than $\text{EMSO}^2(+1, <, \sim)$.*

Proof. The inclusion holds by definition. It is strict because

- the language $L_{a < b}$ cannot be decided by a WDA (Lemma 2) and thus cannot be defined in $\text{EMSO}^2(+1, \sim)$,
- but it can be expressed by the simple formula $\forall x \exists y (a(x) \rightarrow (b(y) \wedge x < y \wedge x \sim y))$.

\square

5 Weak Büchi data automata

In this section we consider automata and logics for *data ω -words*, that is, data words of infinite length. Weak data automata $(\mathcal{A}, \mathcal{C})$ can easily be adapted for data ω -words. The automaton \mathcal{A} is simply interpreted as a letter-to-letter Büchi transducer. A run is accepting if it visits infinitely often a state from F . We refer to the resulting model as weak Büchi data automata (WBDA). We write $\mathcal{L}^\omega(\mathcal{A}, \mathcal{C})$ for the set of data ω -words accepted by $(\mathcal{A}, \mathcal{C})$. The results regarding expressive power of WDA compared with other automata models easily carry over to WBDA.

Data ω -words can be represented by logical structures

$$w = \langle \mathbb{N}, +1, <, \{a(\cdot)\}_{a \in \Sigma}, \sim \rangle, \quad (1)$$

where \mathbb{N} is the set $\{1, 2, \dots\}$ of natural numbers which represent the positions and the other relations are as in the case of data words. For a set $SS \subseteq \{+1, <, \sim\}$ of relation symbols $E_\infty \text{MSO}(SS)$ consists of all formulas of the form

$$\exists_\infty R_1 \dots \exists_\infty R_m \exists S_1 \dots \exists S_\ell \varphi, \quad (2)$$

where $\varphi \in \text{FO}^2(SS)$. Here all relation symbols R_i, S_i are unary. The \exists_∞ are semantically restricted to bind to infinite sets only.

Remark 1. It is folklore that languages (without data) accepted by Büchi automata are precisely languages expressible in formulae of the form:

$$\exists_\infty R_1 \dots \exists_\infty R_m \exists S_1 \dots \exists S_\ell \varphi$$

for some $\varphi \in \text{FO}^2(+1)$. However, we have not found an explicit reference for this result in the literature.

The following theorem is a straightforward generalization of Theorem 5.

Theorem 6. *Weak Büchi data automata and $E_\infty \text{MSO}^2(+1, \sim)$ are equivalent in expressive power.*

Proof. The translation from an automaton to a formula uses one additional relation symbol R which is quantified by an \exists_∞ symbol. The formula φ_{accept} used in the proof of Lemma 3 is then replaced by

$$\forall x (R(x) \rightarrow \bigvee_{q \in F} R_q(x)).$$

For the opposite translation, it can be checked with the help of the Büchi condition that relations quantified with \exists_∞ are indeed infinite. \square

Theorem 7. *The nonemptiness problem for weak Büchi data automata is decidable in 2-NEXPTIME.*

Proof. We show in the following that the nonemptiness problem for WBDA can be polynomially reduced to the nonemptiness problem for WDA. The result then follows from Theorem 2. The approach is a classical one. We show that if the language of a WBDA $(\mathcal{A}, \mathcal{C})$ is non-empty then a finite data string of the form uv can be constructed such that there is a run of \mathcal{A} which loops over v . The “unravelling” uv^ω is then also accepted by the automaton. However, some care is needed to assign data values in a suitable manner.

Let $(\mathcal{A}, \mathcal{C})$ be a WBDA with $\mathcal{A} = \langle \Sigma, \Gamma, Q, q_0, \Delta, F \rangle$. Since we are only interested in whether $\mathcal{L}^\omega(\mathcal{A}, \mathcal{C}) = \emptyset$, we can assume, without loss of generality, that the transitions of \mathcal{A} are all of the form (q, γ, q', γ) . Otherwise, we can replace it by a transducer which reads Γ -strings and guesses, for every position i , a symbol $a_i \in \Sigma$, its profile symbol s_i (and store them in the state) and verifies that its output would be (the actual input symbol) γ_i . Therefore, we consider \mathcal{A} in this proof just as a normal Büchi automaton that gets a Γ -string as input. The constraints are applied to the same string.

We first fix some notation. We refer to the symbols that occur in key constraints of \mathcal{C} as *key symbols*.

A *zone* is a finite data string over Γ in which all positions carry the same data value. An ω -*zone* is an infinite data string over Γ in which all positions carry the same data value. The *zones of a data string* w are the maximal zones of w . An *adorned zone* is a zone together with a pair (q, q') of states of \mathcal{A} . We write $\mathbf{a}\text{-Proj}(z)$ for the triple $(\text{Str}(z), q, q')$ of a zone z that is adorned with the pair (q, q') .

We next define an important notion for this proof, (singular and non-singular) witnesses. We will show that the nonemptiness of $(\mathcal{A}, \mathcal{C})$ boils down to deciding whether such witnesses exist. Singular witnesses correspond to data strings in $\mathcal{L}^\omega(\mathcal{A}, \mathcal{C})$ with an infinite zone whereas non-singular witnesses correspond to data strings with finite zones only.

A *singular witness* for $(\mathcal{A}, \mathcal{C})$ is a data string uv over Γ the following properties.

- $uv \models \mathcal{C}$.
- There is a state $\hat{q} \in F$ and a (partial) run $\rho = \rho_u \rho_v$ of \mathcal{A} on input $\text{Profile}(uv)_\top$ in which the state after reading u and after reading v is \hat{q} . Here, $\text{Profile}(uv)_\top$ denotes the profile string that is obtained from $\text{Profile}(uv)$ by setting the last profile symbol to \top .
- All positions of v and the last zone of u carry the same data value and v does not carry any key symbol.

A *non-singular witness* for $(\mathcal{A}, \mathcal{C})$ is a data string uv over Γ which fulfills the following conditions.

- All zones in uv are of length at most $|Q|(|\Gamma| + 1)$.
- The data value of the last position of u is different from the value of the first position of v .
- There is a state \hat{q} and a (partial) run $\rho = \rho_u \rho_v$ of \mathcal{A} on input uv in which the state after reading u and after reading v is \hat{q} . Furthermore, ρ_v contains some state from F . In the following, each zone z of w is adorned by the pair (q, q') where q is the state of ρ before reading z and q' is the state after reading z .

- The classes of uv can be colored⁴ with the four colors black, yellow, white and blue such that all black, yellow and white classes satisfy⁵ all constraints from \mathcal{C} and furthermore the following conditions hold.
- (black) There are at most $3|Q|^2$ black classes. There are no key symbols in black zones of v . Furthermore, it is not the case that the first zone and the last zone of v are from the same black class.
- (yellow) There are at most $|Q|^2$ yellow classes and they consist of at most $|I|$ zones. All these zones are located in v .
- (white) All zones of the white classes are located in u .
- (blue) For each blue zone z there is a yellow zone z' such that $\mathbf{a}\text{-Proj}(z)=\mathbf{a}\text{-Proj}(z')$.

The proof of decidability of the nonemptiness problem for WBDA now reduces to proving the following three claims.

- (Claim 1) If there exists a witness for $(\mathcal{A}, \mathcal{C})$ then $\mathcal{L}^\omega(\mathcal{A}, \mathcal{C}) \neq \emptyset$.
- (Claim 2) If $\mathcal{L}^\omega(\mathcal{A}, \mathcal{C}) \neq \emptyset$ then there exists a witness for $(\mathcal{A}, \mathcal{C})$.
- (Claim 3) There is a nondeterministic algorithm which constructs, for every WBDA $(\mathcal{A}, \mathcal{C})$, in polynomial time some WDA $(\mathcal{A}', \mathcal{C}')$ such that every possible $(\mathcal{A}', \mathcal{C}')$ accepts only witnesses for $(\mathcal{A}, \mathcal{C})$ and for each witness uv there is a run of the algorithm producing some $(\mathcal{A}', \mathcal{C}')$ that accepts uv .

Therefore, the nonemptiness problem for WBDA can indeed be reduced non-deterministically in polynomial time to the nonemptiness problem for WDA.

Next, we prove Claims 1-3.

We start by proving Claim 1. Let us assume there is a singular witness uv for $(\mathcal{A}, \mathcal{C})$ (where all names are chosen as above). It is easy to see that in this case $\rho_u \rho_v^\omega$ is an accepting run of \mathcal{A} on input uv^ω and that uv^ω satisfies all constraints from \mathcal{C} . It should be noted that in uv^ω all positions in the v^ω -part and some non-empty suffix of u carry the same data value and thus constitute one infinite zone. The repetition of v does not introduce any violations of \mathcal{C} as v does not contain any key symbols.

We next consider the case that uv is a non-singular witness for $(\mathcal{A}, \mathcal{C})$ (where again all names are chosen as above). In principle, we aim again for a data ω -string in $\mathcal{L}^\omega(\mathcal{A}, \mathcal{C})$ that is obtained from uv by repeating v infinitely often. Indeed, by doing so, we obtain a data ω -string whose adorned projection is just $\mathbf{a}\text{-Proj}(u)\mathbf{a}\text{-Proj}(v)^\omega$. However, the data values cannot be the same in every copy of v as otherwise constraints might be violated.

The basic idea for the assignment of data values is as follows. As white zones only appear in u they are not affected and we do not need to adapt them. As the black zones in v do not contain any key symbols, we can leave them unchanged in each of the infinitely many copies of v that constitute w . It only remains to assign data values to the blue zones in u and to the blue and yellow zones in v and the copies of v . To this end, we intuitively use the yellow classes as templates. More precisely, we make sure that for every new class that is constituted by assigning data values, the set of zones corresponds to one of the yellow classes of v , that is, it has the same set of (adorned) zones as that class. This ensures that each new class satisfies \mathcal{C} .

⁴ Each class gets exactly one color. We refer to zones and positions in a black class as black zones and positions, respectively, and likewise for the other colors.

⁵ We do not require that the blue classes satisfy \mathcal{C} .

We now describe the construction of the data ω -string w in more detail.

- Let $w_1 = uv^\omega$. Clearly, $\rho_u \rho_v^\omega$ is an accepting run on $\text{Profile}(uv^\omega)$.
- In the remainder of the construction, only data values are changed, but zone projections and runs remain the same.
- Let w_2 be a copy of w_1 where data values of blue and yellow zones are removed (and thus the black and white zones are just as in u and v). As the last zone of v is not from the same black class as the first zone of v it cannot happen that two black zones with the same data value become adjacent by repeating v .
- Next, we choose an infinite sequence d_1, d_2, \dots of data values that do not occur in black or white zones. We assign data values to the blue and yellow zones in w by repeating the following procedure from left to right. In the i -th application we constitute a new class by assigning the data value d_i to a set of zones that corresponds to some yellow class.
 - We pick the first yellow or blue zone z that does not yet have a data value. We choose a yellow class c of v that contains a zone z' with $\mathbf{a}\text{-Proj}(z) = \mathbf{a}\text{-Proj}(z')$. This is possible as z is either such a zone itself or it was a blue zone in uv and thus such a zone z' exists by the requirements for blue zones. Let z'_1, \dots, z'_k be the other zones of the class c . For each $i \leq k$ we choose some zone z_i of w_2 that has not yet received a data value and fulfills $\mathbf{a}\text{-Proj}(z'_i) = \mathbf{a}\text{-Proj}(z_i)$. We require that the zones z, z_1, \dots, z_k are pairwise not adjacent. As each yellow zone of v is copied infinitely often in w , such zones z_1, \dots, z_k do exist. We assign to z, z_1, \dots, z_k the data value d_i . We note that the new class has exactly the same zone profile as the yellow class c and therefore satisfies all constraints.

We denote the resulting data ω -word by w .

It remains to show that indeed $w \in \mathcal{L}^\omega(\mathcal{A}, \mathcal{C})$. As ρ is an accepting run yielding w it only remains to show that all classes of w satisfy \mathcal{C} .

- As white classes have not been changed at all, they clearly satisfy \mathcal{C} .
- As the black zones in v do not contain any key symbols, repeating them in w does not introduce any violations of key constraints. Otherwise, the repetition does not change the set of occurring symbols for any black class and therefore also the inclusion and denial constraints remain valid.
- Each other class of w has the same set of profiles as some yellow class c of uv . Thus, it satisfies all constraints from \mathcal{C} just as c does.

This concludes the proof of Claim 1.

For the proof of Claim 2 let $\mathcal{L}^\omega(\mathcal{A}, \mathcal{C}) \neq \emptyset$ and let ρ be an accepting run on the data ω -word w with $w \models \mathcal{C}$. We consider two cases, depending on whether the number of zones in w is finite or infinite. We first consider the simpler case, in which the number of zones in w is finite. In this case $w = u'v'$, for some finite data string u' and an infinite data string v' such that

- all positions of v' have the same data value d ,
- there are no key symbols in v' .

The latter can be achieved as there is only a finite number of key symbols in the infinite zone.

As ρ is accepting, some accepting state \hat{q} occurs infinitely often in the run ρ on v' . Let u be the prefix of $u'v'$ until (and including) the position of v' after which \hat{q} occurs for the first time. Let v be the substring of v' from the first to the second occurrence of \hat{q} . Clearly, uv is a singular witness for $(\mathcal{A}, \mathcal{C})$.

Now we turn to the case, where w has an infinite number of zones and therefore all zones are of finite length. The construction of u and v consists of a number of transformation steps of w . We refer to the data string obtained after the i -th transformation step as w_i . In these transformations we intuitively view each w_i as an infinite sequence of zones. We might replace zones by other zones but we never change the sequence of states that ρ takes on the sequence of zone borders. We call the sub-sequence of states that a run ρ takes at zone borders, the *zone sub-run* of ρ .

Thus, our transformations do not change the zone sub-run of our accepting run. Without loss of generality, we assume that whenever \mathcal{A} assumes a state from q after some position p it also assumes an accepting state after the next symbol with profile \perp . This can be accomplished by an easy modification of \mathcal{A} . And now we can be sure that in the zone sub-run some infinite state occurs infinitely often.

The first transformation step transforms w into a data string w_1 in which each zone has length at most $|Q|(|\Gamma| + 1)$. This step is applied to each zone z independently. If $|z| \leq |Q|(|\Gamma| + 1)$, nothing has to be done. Otherwise, we mark a set of positions of z such that the first and last position are marked and, for every symbol γ that occurs, one occurrence is marked. Thus, at most $|\Gamma| + 1$ positions are marked. As $|z| > |Q|(|\Gamma| + 1)$ there must be a sequence of at least $|Q|$ unmarked positions in z . We consider the state of ρ before the first of these positions and after each of them. Clearly, in this sequence of states some state q must occur twice. By removing the data string between the two occurrences of q we obtain a shorter zone z' with the same set of symbols. Furthermore z' can be obtained by a partial run of \mathcal{A} with the same first and last state as for z in ρ , thus the zone sub-run does not change. We note that by removing symbols no key conflicts can be introduced. The repeated application of this process to each zone yields a data word w_1 that is accepted by a run of \mathcal{A} with the same zone sub-run as before and for which $w_1 \models \mathcal{C}$.

We select, for each class c of w_1 , and each symbol a that occurs in c , one zone z of c that contains a . We call these selected zones the *core zones* of c and the other zones *redundant zones*. Clearly, each class has at most $|\Gamma|$ core zones and remaining zones do not contain any key symbols. Thus, if a redundant zone is removed from a class or a copy of a redundant zone is added to a class the validity of constraints is not affected.

In a nutshell, the remaining transformation steps do the following. First of all, we collect all redundant zones in a finite number of classes. These will be the black classes and they are the only classes that might have an infinite number of zones. From the remaining classes we first distinguish those that contain a zone adornment (q, q') that occurs only a finite number of times. These will be the white classes and there are only finitely many of them. The remaining classes consist only of core zones with adornments that occur infinitely often. We single out a polynomial number of such classes, the yellow classes, that cover all “infinite adornments” and in all remaining classes, the blue ones, we replace all zone strings by strings from yellow zones, thereby ensuring that there exists only a polynomial number

of different zone strings outside black and white classes. We now continue the detailed description of the construction.

In the next step, we transform w_1 into a data string w_2 in which at most $3|Q|^2$ classes have redundant zones. Thus, in particular, at most $3|Q|^2$ classes have infinitely many zones. To this end, we proceed as follows, for every pair (q, q') of states of \mathcal{A} . If (q, q') occurs as adornment of any redundant zone of w_1 we pick the (up to) first three classes that contain such zones. We color all these classes black.

Next, we modify all redundant⁶ zones z that are not (yet) in a black class in a left-to-right fashion as follows. Let (q, q') be the adornment of such a zone z . As z is not black there must be three black classes c_1, c_2, c_3 with (q, q') -adorned redundant zones z_1, z_2, z_3 , respectively. We replace z by one of z_1, z_2, z_3 that has a different data value from the zones adjacent to z . Although, this step might change the string projections of zones, the resulting data string still has the same zone sub-run and is therefore still accepted by \mathcal{A} . Furthermore, as only redundant zones were removed from non-black classes, these classes still satisfy \mathcal{C} . And, as in black classes, only copies of redundant zones are added, they also still satisfy \mathcal{C} .

We call a state pair (q, q') *frequent* if it occurs as adornment of infinitely many non-black zones, otherwise *infrequent*. Clearly, there is only a finite number of classes that contain zones with infrequent adornment. We color these classes white.

If all classes are black or white then the construction is finished. Otherwise, the adornment of all zones that are neither black nor white is frequent.

The main goal of the final transformation step is to reduce the number of different string projections that occur in zones that are neither black nor white. We note that this transformation step might cause violations of constraints for blue or yellow classes (but this does not matter as long as we yield a witness). In the following, we choose three positions p_1, p_2, p_3 such that one of them ($p \in \{p_1, p_2\}$) will mark the end of u and such that v will (basically) be the data string between p and one of the other two ($p' \in \{p_2, p_3\}$). We pick three such positions to ensure that the condition on the first and last zone of v holds.

Let \hat{q} be some accepting state that occurs infinitely often in the zone sub-run. Such a state exists by our assumption that the automaton assumes accepting states at the end of a zone if it assumed one inside the zone. Let p_1 be the minimal position in which ρ_2 assumes \hat{q} at the end of a zone and such that all white zones are before p_1 . We next choose, for each frequent pair (q, q') one class $c_{q,q'}$ of w_2 that is neither black nor white, contains a zone with adornment (q, q') and is located after p_1 . Let p_2 be the minimal position in which the zone sub-run assumes \hat{q} and such that all zones of classes $c_{q,q'}$ are before p_2 .

Now, we choose, for each frequent pair (q, q') one class $c'_{q,q'}$ of w_2 that is neither black nor white, contains a zone with adornment (q, q') and is located after p_2 . Finally, we let p_3 be the minimal position in which the zone sub-run assumes \hat{q} and such that all zones of classes $c_{q,q'}$ are before p_3 .

If the first zone after p_1 has a different data value than the last zone before p_2 or at least one of them is not black we set $p = p_1$ and $p' = p_2$. Otherwise, if the first zone after p_2 has a different data value than the last zone before p_3 or at least one of them is not black we set $p = p_2$ and $p' = p_3$. Otherwise, we set $p = p_1$ and $p' = p_3$. In either case, the

⁶ We remind the reader that the term redundant is always relative to a class. We note that the zones z are redundant in their original class and also in their black target class.

first zone after p has a different data value than the last zone before p' or at least one of them is not black.

If $p = p_1$ we color the classes $c_{q,q'}$ yellow, otherwise we color the classes $c'_{q,q'}$ yellow.

In the last transformation step, we color all not yet colored zones in blue and furthermore modify blue zones as follows. Let z be a blue zone with adornment (q, q') . As z is neither white nor black, (q, q') is frequent. Let z' be the (q, q') -adorned zone in $c_{q,q'}$ (or in $c'_{q,q'}$ if $p = p_2$). We keep the data value of z but replace its string projection by $\text{Str}(z')$. This does not affect the zone sub-run, but it might cause a constraint conflict for the blue class (but, as already mentioned, we need not care about this).

Let w_3 be the resulting data string.

Now we define u to be the prefix of w_3 until position p and v to be the data string from (excluding) position p to position⁷ p' .

This construction guarantees that uv is a non-singular witness for $(\mathcal{A}, \mathcal{C})$. This completes the proof of Claim 2.

Finally, we prove Claim 3. As we assume that \mathcal{A} copies its input string to the output (and thus basically is an automaton) the reduction also constructs WDA with this property. We first show how to compute a WDA $(\mathcal{A}', \mathcal{C}')$ for singular witnesses for $(\mathcal{A}, \mathcal{C})$. The algorithm first guesses a symbol a_0 that occurs in the infinite class.

\mathcal{A}' has input alphabet $\Gamma \times \{0, 1, 2\}$. The symbols of the form $(a, 1)$ and $(a, 2)$ are used for the class of v and the others for the remaining classes. \mathcal{A}' simulates \mathcal{A} when it reads only the Γ -part. Furthermore, it guesses a position p (intuitively: the border between u and v) that has a state \hat{q} from F and verifies that the final state is \hat{q} as well. For the simulation of the last step \mathcal{A}' behaves as if the final symbol carried a \top -symbol in its profile part. \mathcal{A}' further checks that all symbols from $\Gamma \times \{0\}$ occur before position $p - 1$ and that no symbols $(a, 1)$ or $(a, 2)$, where a is a key symbol of \mathcal{C} occur after p . It furthermore checks that there is only one occurrence of $(a_0, 2)$ and no other $(a, 2)$. All constraints from \mathcal{C} are reproduced in \mathcal{C}_1 separately, for symbols of the form $(a, 0)$ and $(a, 1)$. Furthermore, they ensure that in each class either only $(a, 0)$ -symbols occur or only $(a, 1)$ -symbols and $(a_0, 2)$ (by denial constraints). Finally, there is an inclusion constraint $V((a, 1)) \subseteq V((a_0, 2))$, for every $a \in \Gamma$ and a key constraint for $(a_0, 2)$, making sure that there is only one class with symbols from $\Gamma \times \{1\}$.

It remains to show how to (non-deterministically) compute a WDA $(\mathcal{A}', \mathcal{C}')$ for non-singular witnesses for $(\mathcal{A}, \mathcal{C})$. The basic idea is that the algorithm first guesses the adorned zones that are used in yellow classes in the order in which they appear in the witness. Furthermore, for each black class, it guesses some symbol a occurring in that class and it guesses the order in which these symbols occur. These symbols are colored with black' instead of black. We use these non-deterministic guesses in the reduction as otherwise, \mathcal{A}' would need to handle, e.g., all possible orders in which the yellow zones appear. This would result in an exponential blow-up.

⁷ Formally, the position p' might have been modified during the last transformation step. However, we refer by p' to the last position of the zone that corresponds to the zone that ended in p' in w_2 .

\mathcal{A}' uses the alphabet $\Gamma \times \{\text{black}, \text{black}', \text{yellow}, \text{white}, \text{blue}\} \times \{0, \dots, 3|Q|^2\}$. It reads colored symbols and always simulates \mathcal{A} on the uncolored projection. It guesses a position p and a state \hat{q} and checks that

- after position p the run has state \hat{q} and likewise at the end;
- between position p and the end, the run assumes some accepting state;
- white symbols only occur before p ;
- the yellow zones all appear after p and they exactly correspond to the adorned zones that were guessed before;
- each blue zone has the same adorned projection as some yellow zone;
- white, yellow and blue symbols carry a 0 in their last component;
- black symbols carry a non-zero number in their last component;
- there are no black key symbols after p ;
- it is not the case that the first zone after p is in the same class as the last zone and that they are both black;
- each expected black' symbol occurs exactly once and they occur in the expected order.

In \mathcal{C}' the constraints of \mathcal{C} are reproduced for the black, yellow and white class. Some constraints are added that ensure that in each black class a black'-symbol occurs, similarly as for the singular case. Furthermore, in black classes, all symbols have the same number in their last component. In this way, it is ensured that for each $i \in \{1, \dots, 3|Q|^2\}$, there is at most one black class.

Clearly, \mathcal{A}' and \mathcal{C}' can be computed in polynomial time. The computation is deterministic, once the above mentioned values are guessed.

This completes the proof of Claim 3 and thus the proof of the theorem. \square

6 Conclusion

We conclude this paper with two open problems for future directions. An obvious open problem is the exact complexity of the nonemptiness problem for weak data automata. The current 2-NEXPTIME yields a 3-NEXPTIME upper bound for the satisfiability problem for $\text{EMSO}^2(+1, \sim)$. However, as it is known that this problem can be solved in 2-NEXPTIME [18], some room for improvement is left.

Another interesting question is how our results can be applied to temporal logics. In [10], a restriction of LTL with one register, *simple* LTL, was considered with the same expressive power as some two variable logic. We conjecture that there is a correspondence between our logics and the restriction of simple LTL to the operators X , X^{-1} and an operator that allows navigation to some other position.

References

1. Henrik Björklund and Thomas Schwentick. On notions of regularity for data languages. *Theor. Comput. Sci.*, 411(4-5):702–715, 2010.
2. Luc Boasson. Some applications of CFL's over infinite alphabets. In *Theoretical Computer Science*, pages 146–151, 1981.

3. Mikolaj Bojanczyk, Anca Muscholl, Thomas Schwentick, and Luc Segoufin. Two-variable logic on data trees and XML reasoning. *J. ACM*, 56(3), 2009.
4. Mikolaj Bojanczyk, Anca Muscholl, Thomas Schwentick, Luc Segoufin, and Claire David. Two-variable logic on words with data. In *LICS*, pages 7–16, 2006.
5. Benedikt Bollig. An automaton over data words that captures EMSO logic. *CoRR*, abs/1101.4475, 2011.
6. J. R. Büchi. Weak second-order arithmetic and finite automata. *Z. Math. Logik Grundl. Math.*, 6:66–92, 1960.
7. Edward Y. C. Cheng and Michael Kaminski. Context-free languages over infinite alphabets. *Acta Inf.*, 35(3):245–267, 1998.
8. Claire David, Leonid Libkin, and Tony Tan. On the satisfiability of two-variable logic over data words. In *LPAR (Yogyakarta)*, pages 248–262, 2010.
9. Stéphane Demri, Deepak D’Souza, and Régis Gascon. A decidable temporal logic of repeating values. In *LFCS*, pages 180–194, 2007.
10. Stéphane Demri and Ranko Lazic. LTL with the freeze quantifier and register automata. *ACM Trans. Comput. Log.*, 10(3), 2009.
11. Calvin C. Elgot. Decision problems of finite automata design and related arithmetics. *Transactions of The American Mathematical Society*, 98:21–21, 1961.
12. Jay L. Gischer. Shuffle languages, Petri nets, and context-sensitive grammars. *Commun. ACM*, 24(9):597–605, 1981.
13. Erich Grädel and Martin Otto. On logics with two variables. *Theor. Comput. Sci.*, 224(1-2):73–113, 1999.
14. Michael Kaminski and Nissim Francez. Finite-memory automata. *Theor. Comput. Sci.*, 134(2):329–363, 1994.
15. Michael Kaminski and Tony Tan. Regular expressions for languages over infinite alphabets. *Fundam. Inform.*, 69(3):301–318, 2006.
16. Ranko Lazic. Safety alternating automata on data words. *ACM Trans. Comput. Log.*, 12(2):10, 2011.
17. Frank Neven, Thomas Schwentick, and Victor Vianu. Finite state machines for strings over infinite alphabets. *ACM Trans. Comput. Log.*, 5(3):403–435, 2004.
18. Matthias Niewerth and Thomas Schwentick. Two-variable logic and key constraints on data words. In *ICDT*, pages 138–149, 2011.
19. Friedrich Otto. Classes of regular and context-free languages over countably infinite alphabets. *Discrete Applied Mathematics*, 12(1):41 – 56, 1985.
20. Wolfgang Thomas. Languages, automata, and logic. In G. Rozenberg and A. Salomaa, editors, *Handbook of Formal Languages, Vol. III*, pages 389–455. Springer, New York, 1997.
21. Boris Trakhtenbrot. Finite automata and logic of monadic predicates. *Doklady Akademii Nauk SSSR*, 140:326–329, 1961.